

Cluster-wise peak detection and filtering based on spatial distribution to efficiently mine mass spectrometry imaging data

Jonatan Eriksson, Melinda Rezeli, Max Hefner, György Marko-Varga, and Péter Horvatovich

Anal. Chem., **Just Accepted Manuscript** • DOI: 10.1021/acs.analchem.9b02637 • Publication Date (Web): 12 Aug 2019

Downloaded from pubs.acs.org on August 13, 2019

Just Accepted

“Just Accepted” manuscripts have been peer-reviewed and accepted for publication. They are posted online prior to technical editing, formatting for publication and author proofing. The American Chemical Society provides “Just Accepted” as a service to the research community to expedite the dissemination of scientific material as soon as possible after acceptance. “Just Accepted” manuscripts appear in full in PDF format accompanied by an HTML abstract. “Just Accepted” manuscripts have been fully peer reviewed, but should not be considered the official version of record. They are citable by the Digital Object Identifier (DOI®). “Just Accepted” is an optional service offered to authors. Therefore, the “Just Accepted” Web site may not include all articles that will be published in the journal. After a manuscript is technically edited and formatted, it will be removed from the “Just Accepted” Web site and published as an ASAP article. Note that technical editing may introduce minor changes to the manuscript text and/or graphics which could affect content, and all legal disclaimers and ethical guidelines that apply to the journal pertain. ACS cannot be held responsible for errors or consequences arising from the use of information contained in these “Just Accepted” manuscripts.

Cluster-wise peak detection and filtering based on spatial distribution to efficiently mine mass spectrometry imaging data

Jonatan O. Eriksson,[†] Melinda Rezeli,[†] Max Hefner,[†] Gyorgy Marko-Varga,[†] and
Peter Horvatovich^{*,‡,†}

[†]*Lund University, Department of Biomedical Engineering, Lund*

[‡]*University of Groningen, Department of Analytical Biochemistry, Groningen Research
Institute of Pharmacy, Antonius Deusinglaan 1, 9713 AV, Groningen, the Netherlands*

E-mail: p.l.horvatovich@rug.nl

Abstract

Mass spectrometry imaging (MSI) has the potential to reveal the localization of thousands of biomolecules such as metabolites and lipids in tissue sections. The increase in both mass and spatial resolution of today's instruments brings on considerable challenges in terms of data processing; accurately extracting meaningful signals from the large data sets generated by MSI without losing information that could be clinically relevant is one of the most fundamental tasks of analysis software. Ion images of the biomolecules are generated by visualizing their intensities in 2-D space using mass spectra collected across the tissue section. The intensities are often calculated by summing each compound's signal between predefined set of borders (bins) in the m/z dimension. This approach, however, can result in mixed signals from different compounds in the same bin or splitting the signal from one compound between two adjacent bins leading to low quality ion images. To remedy this problem, we propose

1
2
3 a novel data processing approach. Our approach consists of a sensitive peak detection
4 method able to discover both faint and localized signals by utilizing cluster-wise kernel
5 density estimates (KDEs) of peak distributions. We show that our method can recall
6 more ground-truth molecules, molecule fragments, and isotopes than existing methods
7 based on binning. Furthermore, it automatically detects previously reported molecular
8 ions of lipids, including those close in m/z , in an experimental data set.
9
10
11
12
13
14
15
16

17 Introduction

18
19
20 Mass spectrometry imaging (MSI) is a technique often used to study the localization of
21 known and unknown biomolecules such as lipids, metabolites, or peptides in tissue. Today's
22 instruments can scan samples with both high spatial and mass spectral resolution and, con-
23 sequently, generate massive data sets that require highly efficient and accurate processing.
24 Thus, one of the key component of MSI data processing is data-reduction, which typi-
25 cally involves detection and extraction of signals originating from tissue or drug compounds
26 while discarding noise.^{1,2} The peaks of each spectrum are mapped onto a common reference
27 and by visualizing the intensities of individual peaks as images the spatial distribution of
28 biomolecules can be revealed. The reference spectrum is generated by detecting peaks which
29 are common to multiple spectra. Accurate peak detection facilitates the isolation of signals
30 from individual compounds which is necessary to obtain high quality images.
31
32
33
34
35
36
37
38
39
40
41

42 Many existing MSI software, such as Cardinal³ and MALDIquant⁴ detect isotopic peaks
43 of compounds on a data set mean spectrum and subsequently rank them based on the
44 frequency of their presence in ion image pixels. This method is fast and produces concise peak
45 lists but has limited performance for low-intensity peaks and those localized to small regions
46 in the analyzed tissue section.¹ Many tools generate ion images by binning around each peak
47 of interest; the intensity value for each pixel is calculated by summing ion intensities between
48 predefined m/z borders (bins). When doing this, however, it is crucial to use narrow bins to
49 avoid mixing signals from multiple compounds in one image and to ensure that the mass of
50
51
52
53
54
55
56
57
58
59
60

1
2
3 the peak around which binning is performed is accurate.
4

5 Suits et al.⁵ showed that *slicing* the entire m/z range into ion images of fixed mass
6 widths enables MSI practitioners to explore MSI data sets in a hypothesis-free manner. This
7 approach sets no threshold on either peak intensity or presence in a minimum number of
8 pixels and is thus not biased toward large or high intensity molecules in the tissue. Choosing
9 bin width is a specificity-sensitivity trade off. A small bin width results in higher sensitivity
10 but increases the risk of peak splitting and a higher number of empty or non-informative
11 ion images. Larger bin widths on the other hand result in fewer non-informative images
12 but are unable to discriminate between compounds that are close in mass, resulting in
13 ion images containing signals from multiple compounds. Unfortunately, even when using
14 relatively large bin widths, slicing leads to impractically large sets of ion-images unless the
15 experimentalist is guided by known ion masses. However, previous studies have demonstrated
16 that incorporating information about the ion-images' spatial structure in MSI data analysis
17 pipelines is an effective way to automatically separate high and low quality images in these
18 large image sets.⁶⁻⁹
19
20
21
22
23
24
25
26
27
28
29
30
31
32

33 In this paper, we present a peak detection method that enables automatic detection
34 of faint and localized signals as well as high intensity and/or abundant signals. We show
35 that our peak detection can serve as a part of a MSI data analysis pipeline that is both
36 sensitive and specific by combining it with established methods that filter peaks based on
37 their spatial arrangement. A sensitive peak detection algorithm is not only essential for
38 exploratory analysis, but also for discovering molecules spatially co-localized with those
39 expected to be present, e.g. drug compounds and metabolites. This is highly relevant in
40 both scientific and clinical settings where drug-tissue interaction and tissue composition are
41 often investigated. To assess and compare the performance of our method to existing MSI
42 data processing tools, we used a rat liver section spiked with several drugs, most of which
43 are anticancer drugs, where the masses of the spiked drugs are used as ground-truth. Using
44 this data set, we show that we are able to detect drug peaks as well as fragment and isotopic
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

1
2
3 peaks, including those that are close in m/z to more intensive and/or abundant peaks. We
4
5 also used the MSI data set from a mouse bladder section originally presented by Römpp
6
7 et al.¹⁰ to further assess our method.
8
9

11 Materials and methods

15 Drug compounds and matrix composition

18 For the MALDI-MSI experiment, we selected 12 different drugs. The drugs were purchased
19
20 from the LC Laboratories (Woburn, MA, USA; CAS numbers: dabrafenib: 1195765-45-7,
21
22 dasatinib: 302962-49-8, erlotinib: 183321-74-6, gefitinib: 184475-35-2, imatinib: 152459-
23
24 95-5, lapatinib: 388082-78-8, pazopanib: 444731-52-6, sorafenib: 284461-73-0, sunitinib:
25
26 557795-19-4, trametinib: 871700-17-3, vatalanib: 212141-54-3) and from SelleckChem (Mu-
27
28 nich, Germany; CAS numbers: ipratropium: 60205-81-4) with >99% purity and were dis-
29
30 solved in methanol (MeOH, (ChromasolvTM Plus for HPLC) (Sigma-Aldrich, Steinheim,
31
32 Germany) at 10 mg/mL concentration. These stock solutions were further diluted with
33
34 50% MeOH and 5 mixtures were generated, each containing 4 different drug compounds.
35
36 Supplementary table 1 summarizes the composition of the 5 drug mixtures. 5 mg/mL α -
37
38 cyano-4-hydroxycinnamic acid (CHCA, Sigma-Aldrich) dissolved in 50% MeOH containing
39
40 0.1% trifluoroacetic acid (TFA, Sigma-Aldrich, Steinheim, Germany) was used as matrix
41
42 solution.
43
44

45 Sample preparation

48 For MALDI-MSI, a 10 μm section was cut from frozen rat liver tissue using a cryotome
49
50 and placed on a glass slide. Then 0.3 μL from each drug mixture was pipetted on the
51
52 tissue section at predefined positions. After drying of the tissue, CHCA matrix solution
53
54 was deposited on the tissue surface by an automated pneumatic sprayer (TM-Sprayer, HTX
55
56 Technologies). The nozzle distance was 46 mm , and the spraying temperature was set to
57
58

1
2
3 35 °C, the matrix was sprayed (19 passes) over the tissue section at a linear velocity of 750
4 *mm/min* with a flow rate set to 0.1 *ml/min* and a nitrogen pressure set at 10 *psi*. After
5
6
7 each pass, a drying time of 30 *s* was set on the spraying machine to give time for the sample
8
9 to dry completely before the next pass. The frozen rat liver tissue was provided by prof.
10
11 Roland Andersson (Dept. Clinical Sciences Lund (Surgery), Skane University Hospital, Lund
12
13 University). Animals were housed and bred according to regulations for the protection of
14
15 laboratory animals.
16
17
18

19 MALDI MSI

20
21
22 MSI data was collected by sampling the tissue section with 50 μm raster arrays without
23
24 laser movement within each measuring position. The dimensions of the measured liver tissue
25
26 section was approximately 0.9 by 1.2 *cm* in x, y sampling coordinates. A total of 23,823
27
28 sampling positions ($x = 247$, $y = 181$) were collected. Full mass spectra were collected by
29
30 using a MALDI LTQ Orbitrap XL mass spectrometer (Thermo Fisher Scientific, Bremen,
31
32 Germany), equipped with a 60 Hz 337 nm nitrogen pulse laser (LTB Lasertechnik Berlin,
33
34 Berlin, Germany). This instrument was operated at 60 000 resolution (at m/z 400) collecting
35
36 spectral data in the mass range of 150-1000 m/z in profile mode generated by 20 laser shots
37
38 at 10 μJ with automatic gain control switched off. Data were acquired using Xcalibur v
39
40 2.0.7. software (Thermo Fisher Scientific, San Jose, CA). The MSI raw data contains mass
41
42 spectra from all measurement points together with their x, y coordinates.
43

44 The Thermo Scientific raw files were first converted to *mzML* using *msconvert* and then to
45
46 *imzML*¹¹ format using *imzmlConverter*. Finally, the *imzML* data was loaded into MATLAB
47
48 and analyzed with custom scripts. The mouse bladder data set with PXD001283 ID was
49
50 downloaded from ProteomeXchange in *imzML* format.
51
52
53
54
55
56
57
58
59
60

Peak picking

We propose a two step peak picking scheme; in the first step, candidate peaks are detected on clusters of peak m/z values from all spectra, and in the second, the spatial distribution of the candidate peaks is evaluated and we select those that display a coherent structure. For the first step, we have devised a novel method that relies on cluster-wise kernel density estimates (KDEs) of spectral peaks. KDEs are smooth histograms and we use them to estimate the distribution of the peak m/z values within clusters along the m/z axis. The level of smoothness is adapted to each cluster independently. Candidates of data set peaks are then detected as local maxima on the resulting KDE curves. For the second step, we use two established ways to automatically estimate the quality of the images corresponding to peaks obtained in the first step as a means to filter out non-informative peaks. **Figure 1** summarizes all parts of our peak picking scheme.

Peak Detection

Firstly, we collect the peak masses from every spectrum in one list, mz_{all} , which is then sorted in ascending order. Centroided spectra are taken as input and peaks with heights below a very low intensity threshold are discarded to reduce the impact of background noise. Consequently, mz_{all} will contain most peak masses from the data set. Depending on data set size and RAM availability mz_{all} is processed either in segments or in its entirety. Secondly, peak clusters in the m/z dimension are identified using a one-dimensional directional graph. If the distance between an m/z value, m_i , and the next, m_{i+1} , is smaller than d_c , an edge connecting the two is added to the graph. The connected components in the resulting graph represent the m/z clusters. We let d_c increase with m/z to account for the peak broadening described by the known theoretical relationship between peak width (at half maximum) and m/z : $d_c = f(m/z)$ where f depends on instrument type.¹² Suits et al.¹³ summarized the relationship between peak width and instrument type. To reduce processing time we discard clusters containing fewer than a minimum number of peaks. The threshold should

1
2
3 be set sufficiently low to retain peaks representing meaningful anatomical structures in the
4 tissue and is therefore dependent on the spatial resolution of the experiment. Finally, to test
5 whether a cluster contains one or more peaks, a KDE is fitted to the distribution of m/z
6 values within the cluster. The kernel bandwidth is optimized for each cluster individually
7 using the normal optimal smoothing method described by Bowman and Azzalini¹⁴. Peaks
8 are detected on the KDE curve in an iterative fashion: first the local maxima are detected
9 and added together with their corresponding heights to a cluster-specific peak list, p_{kde} . The
10 m/z corresponding to the highest peak in this list, mz_{max} , is added to the global peak list,
11 mz_{ref} , and all surrounding peaks in p_{kde} , that fall within d_{kde} including mz_{max} , are removed.
12 This step is repeated until p_{kde} is empty. d_{kde} is proportional to the expected peak width of
13 the instrument in the same manner as d_c . The ion images are then generated by aligning
14 each centroided spectrum to the resulting reference spectrum mz_{ref} , using a nearest neighbor
15 method with maximum drift threshold dependent on the expected theoretical peak width
16 (at half maximum), similarly to the threshold used when generating edges between peaks in
17 the clustering step.
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33

34 Peak Selection

35
36
37 Although our method is more directed than slicing the spectra across the m/z range (since it
38 only considers a selection of the m/z regions) it still generates many peaks representing noise
39 in addition to those correlated with actual tissue structures, making it is essential to separate
40 the former from the latter. We use the spatial chaos⁸ (SC) and the principal component
41 analysis (PCA) based variance explained¹⁵ (VE) measures to automatically estimate the level
42 of structure in the ion images. The spatial chaos counts the number of connected objects in
43 an ion image. More structured ion images are expected to have fewer disconnected (separate)
44 objects than unstructured ones. The VE measure is the percentage of total variance explained
45 by the first pair of singular vectors of each ion image. This corresponds to how much of the
46 variation in intensity along one axis of the image is explained by the intensities along the
47
48
49
50
51
52
53
54
55
56
57
58
59
60

1
2
3 other. The first principal component inherently explains the most variance and, thus, if it
4 explains very little, so will all others. In structured images there is typically an intensity
5 relationship between the axes and therefore their VE is expected to be higher than that
6 of images with randomly distributed intensities, i.e. unstructured images, in which this
7 relationship is unlikely to exist.
8
9
10
11
12
13
14

15 Results and discussion

16
17
18 Two data sets were used to assess the performance of our novel MSI data pre-processing
19 algorithm based on cluster-wise peak detection. The first MALDI-MSI data set (referred
20 to as the "spiked data set") was generated by spiking a rat liver section with 5 mixtures
21 of 4 ground-truth drugs (12 different compounds in total) in various concentrations. These
22 mixtures were spotted on a rat liver tissue section at five different locations in circular
23 areas of the same size (**figure S1**) and, after matrix deposition, the whole tissue section
24 was analyzed by MALDI-MSI using 50 μm spatial resolution. The concentrations of the
25 drug compounds covered an intensity range of 3 orders of magnitude between trametinib
26 (1.70×10^4) and ipratropium (1.49×10^7). Furthermore, some of the ground-truth drugs
27 such as erlotinib and dasatinib, were spotted in multiple spots at different concentrations.
28 The second data set, originally from Römpp et al.¹⁰, comes from a mouse bladder section
29 and was downloaded from ProteomeXchange (XD001283). This MSI data set was generated
30 by a LTQ Orbitrap instrument with an ion source built in-house used to scan the mouse
31 bladder section with 10 μm spatial resolution. The authors of this study presented the ion
32 images of 11 compounds. These images were generated with a narrow bin width of 0.01 Da.
33 For this data set, we use the mass of these compounds as ground truth, i.e. peaks known to
34 be present.
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

Recall of known compounds

We applied Cardinal, MALDIquant, slicing the spectra into 0.05 Da bins, and our cluster-wise peak detection method to the spiked data set to compare their ability to recall compounds. The difference between the known mass of each ground-truth drug and the mass of the closest detected peak is used as the measure of accuracy for Cardinal and our method. The ion images corresponding to the monoisotopic peak of the ground-truth drugs were manually evaluated to confirm that a compound had been correctly found. Firstly, we ran Cardinal and detected 4,751 peaks; we did not filter out those with too low pixel frequency. The corresponding ion images were generated by binning around each peak. Eight of the 12 compounds were detected with a mass deviation ranging between 4.23 and 198.85 ppm (mean 83.983 ppm). **Figure S2** shows the ion images of the drug compounds generated by cardinal. The ion images of erlotinib (394.176 Da) and gefitinib (447.160 Da) are contaminated with signal from other compounds while sunitinib (399.220 Da), imatinib (494.267 Da), and trametinib (616.086 Da) are completely missed. Secondly, we used MALDIquant to compute a mean spectrum on which we detected 521 peaks. Only the peak from the drug with the highest measured intensity, ipratropium, was found with a mass deviation of 4.7145 ppm. The ion image corresponding to the monoisotopic peak of ipratropium indicates that this compound has diffused from the spotting location and because of this covers a significantly larger region of the tissue than the other compounds; this might contribute to its presence in the mean spectrum which favors signals that have high intensity and/or pixel frequency. Thirdly, we sliced the spectra with a bin width of 0.05 Da across the 150 - 1,000 m/z range resulting in 17,000 slices. To assess the sensitivity of the slicing approach we manually examined the ion images corresponding to the slices containing the m/z of the spiked-in drug compounds (**figure S3**). The signal from trametinib (616.086) is missed and those from erlotinib (394.176 Da) and imatinib (494.267 Da) are mixed with others resulting in contaminated ion images. Finally, when applying our method we identified 3,148 m/z clusters in the data set peak list and on the KDEs of these we detected 6,088 peaks. We

1
2
3 used a value of 0.2 times the theoretical peak width at half maximum for d_c , the parameter
4 controlling the maximum distance between connected points that form the m/z clusters. De-
5 creasing or increasing d_c between 0.1 and 0.5 results in a higher or lower number of clusters,
6 respectively, but ultimately has little impact on the final peak list. All of the 12 spiked-in
7 compounds are detected with mass deviations ranging between 1.00 and 4.29 ppm (mean
8 2.598 ppm). **Figure S4** shows the ion images corresponding to the monoisotopic peaks of
9 the drug compounds generated by our method. The signal from trametinib is weak but
10 detected nevertheless; it had the lowest measured intensity which can explain its absence
11 in some of the spectra. Generally, the quality of images generated with our approach is
12 higher than that of the images generated with Cardinal or by slicing. The drug signals are
13 clearly visible against the background and there is no contamination with signals from other
14 compounds, background, or matrix. **Table S1** shows the mass deviations of the detected
15 peaks corresponding to the spiked-in drugs obtained with our algorithm and Cardinal. The
16 corresponding ion images are shown in **figure S4** and **figure S2**, respectively.

17
18
19
20
21
22
23
24
25
26
27
28
29
30
31 An example of a cluster containing densely located molecule signals is that containing
32 erlotinib (394.176 Da) (**figure 2a**). There are four distinctive signals within this relatively
33 narrow m/z window (0.04 Da) at 394.161, 394.166, 394.172, and 394.176 m/z with inter-
34 peak distances of 13, 15, and 10 ppm. The peak at 394.161 m/z is tissue-derived while those
35 at 394.166 m/z and 394.172 come from a fragment molecule of imatinib and the matrix,
36 respectively. Using our method we are able separate the four peaks and generate a clean
37 image for each of them. **Figure 2b - 2e** show the ion images related to these peaks. If the
38 spectra are binned between 394.150 and 394.200 m/z instead, the signals from three of the
39 four compounds appear in the same ion image, i.e. they are incorrectly combined into one
40 ion-image while that from the peak at 394.172 m/z is invisible (**figure 2f**) due to its low
41 intensity compared to the other three. We found that a value between 0.25 – 0.5 times the
42 theoretical peak width at half maximum is a good choice for d_{kde} , the parameter controlling
43 the minimum distance between two adjacent peaks on the KDE curve. Using a higher value
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

1
2
3 results in fewer noise peaks, however, we lose true peaks, e.g. those from imatinib and
4 erlotinib. Because of this, we recommend using a small d_{kde} to delay filtering out noise peaks
5 until after alignment by using one of the spatial distribution based peak selection methods.
6
7 The kernel bandwidth used when generating the cluster KDEs is optimized for each cluster
8 individually to account for the variability in peak density. This parameter determines the
9 level of smoothing when estimating the distribution of the peak masses within the clusters.
10
11 Similarly to d_{kde} , using a higher bandwidth results in less noisy data, however, may lead to
12 losing true peaks or mixing signals from multiple compounds.
13
14
15
16
17
18

19 We also applied our cluster-based peak detection method to the high spatial resolution
20 mouse bladder data set. In this data set we detected 1,702 m/z clusters and 6,482 peaks. We
21 then filtered out peaks which were present in fewer than 200 of the 33,000 spectra resulting
22 in a final list of 1,024 data set peaks. The original paper reported eleven ion images that were
23 manually generated by binning around peaks with known m/z using a very narrow bin width
24 of 0.01 Da. All peaks corresponding to these ion images are found by our peak detection
25 method in an unsupervised fashion, including the two densely located peaks at 770.5097
26 and 770.5698 m/z originating from the K⁺ adduct of PC(32:1) [phosphatidylcholin] and an
27 isotope of the K⁺ adduct of SM(36:1), [sphingosylphosphorylcholine] respectively (**figure**
28 **3**). **Figure S5** shows the ion images related to the eleven detected peaks.
29
30
31
32
33
34
35
36
37
38
39
40

41 Peak selection

42
43 As previously mentioned, we find more than 6,000 peaks in the rat liver data set with
44 our cluster-based peak detection, resulting in an equal number of ion-images. Manually
45 evaluating each image is impractically slow, but by computing the spatial chaos (SC) and
46 the variance explained (VE) for all ion images, including those of the compounds known to
47 be present, we can estimate how much we can reduce the number of images without losing
48 relevant information. For each data set, we took the VE and SC scores of the ion images
49 corresponding to the known compounds and used their mean scores minus two standard
50
51
52
53
54
55
56
57
58
59
60

1
2
3 deviations as low-end thresholds. The number of peaks whose images had scores above these
4 thresholds indicates how many of the detected peaks should be kept and how many can be
5 rejected as noise. In the spiked data set this filtering resulted in a final list of 843 and 2
6
7 170 peaks when we filtered based on VE and SC scores, respectively. The numbers of peaks
8
9 obtained for the mouse bladder data set are 418 and 288 for VE and SC, respectively. The
10
11 number of ion images whose VE or SC score is above various thresholds is shown in **figure**
12
13 **4**. The number of peaks can potentially be further reduced if off-tissue regions are available;
14
15 biologically irrelevant peaks, such as those coming from solvents or the matrix, can be filtered
16
17 out since their signal often is stronger in these regions.¹⁵
18
19
20

21 Despite its simplicity, the VE score proved to be very effective in ranking the quality of
22
23 the ion images generated from both the spiked and mouse bladder data sets. Specifically,
24
25 VE favors images which have intensities localized to small regions, e.g. all of the spiked-in
26
27 compounds in the spiked data set and heme b, M^+ at $m/z = 616$ (**figure S5c**) in the
28
29 mouse bladder data set. In contrast, ion images with high levels of structure across the
30
31 entire scanned region tend to be rewarded with the highest SC scores making it suitable
32
33 as general measure of image quality but less effective than the VE score in identifying ion
34
35 images with localized structured intensity patterns. The two scores appeared to be partially
36
37 complementary to each other; the Pearson correlation between the VE and SC scores in
38
39 the spiked and mouse bladder data sets were 0.6158 and 0.4821, respectively. **Table 1** and
40
41 **table 2** show the VE and SC scores of the ion images corresponding to the ground truth
42
43 compounds in the spiked and mouse bladder data sets, respectively.
44
45
46

47 **Detection of fragments and isotopes**

48

49 MALDI-MSI is an important tool often used to investigate the distribution of drugs and
50
51 drug metabolites in tissue during pharmaceutical research, and obtaining comprehensive list
52
53 of interacting molecules is crucial during their development. To this end, we further assessed
54
55 the performance of our peak detection method by searching for molecules co-localized with
56
57
58
59
60

1
2
3 the drugs in the spiked data set. Co-localization analysis can be performed by computing
4 the Pearson correlation coefficient between the ion image of a peak of interest and all other
5 images.^{5,16,17} For each drug compound we computed the correlation coefficient between the
6 ion image corresponding to its monoisotopic peak and every ion image from the full image
7 sets generated using the peaks found with our cluster-wise peak detection method and that
8 generated by slicing, without performing peak filtering based on spatial distribution. We
9 manually assessed images whose correlation coefficient was ≥ 0.5 to search for candidate
10 fragments and isotopes with spatial intensity distributions matching those of the drugs. The
11 m/z of the matching images and existing knowledge about the theoretical fragmentation
12 pattern of the drugs were then used to identify the fragments. This resulted in the identifi-
13 cation of 46 isotopes and fragments in the ion image set generated by our method and 32 in
14 the set generated by slicing. We gain an additional 14 fragments and isotopes when using
15 our peak detection approach compared to when slicing the spectra with a bin width of 0.05
16 Da.

17
18
19
20
21
22
23
24
25
26
27
28
29
30
31 The correlation analysis result of dasatinib is shown in **figure 5**. In total, 12 ion images
32 have a correlation coefficient ≥ 0.5 . The nine most correlated images (≥ 0.75) consist of
33 three isotopes of dasatinib with an m/z of 489.165, 490.159, and 491.162, and six fragments
34 with an m/z of 319.133, 387.078, 401.094, 402.097, 403.091, and 427.110. The fragments'
35 and isotopes' ion images show minimal signal mixing with other compounds as shown in
36 **figure 5**. The remaining three consist of another fragment of dasatinib with with an m/z of
37 429.106 and a correlation coefficient of 0.5422 and two ion images related to sorafenib. The
38 identified fragments and results of the correlation analysis are presented in **supplementary**
39 **table 4** and **figures S6-S16**. We also assessed the most anticorrelated images to investigate
40 whether there was evidence of ion suppression from any of the ground-truth drugs. However,
41 no images uniquely anticorrelated to any one of the spiking spots were found. Instead, these
42 images were anticorrelated to all spiking spots simultaneously indicating that they are the
43 result of washing or ion suppression from the solvent used in the drug mixtures.

Conclusions

In this paper we have presented an efficient peak picking approach combining a novel peak detection algorithm with filtering based on spatial information to automatically identify ion images corresponding to isotopic peaks of both endogenous and drug compounds in high-resolution MSI data sets. It should be noted that these data sets were generated using high-resolution Orbitrap MSI, which is low-pass filtered during acquisition by default. Applying our method to noisier data such as that generated by QTOF MSI would require additional pre-processing such as baseline removal and smoothing. Our KDE cluster-wise peak detection algorithm enables us to find low intensity and localized peaks with minimal contamination from other peaks close in m/z , resulting in high ion image quality. We believe that implementing our MSI pre-processing algorithm in an interactive tool would be valuable to experimentalists who aim to identify *a priori* unknown endogenous compounds, reveal drug distributions in tissue, or find compounds that spatially correlate to known ones. Such a tool could help users gain deeper insight into the effect of drugs in tissue and considerably reduce the number of ion images that have to be examined manually.

Table 1: The VE and SC scores of the ion images corresponding to the spiked-in drug compound in the spiked data set and their corresponding rank among the 4 771 ion images that remain after removing those with fewer than 400 nonzero pixels.

Compound	Mass	VE	Percentile	Rank (VE)	SC	Percentile	Rank (SC)
Ipratropium	332.223	0.5997	99.43	27	0.9997	99.94	3
Vatalanib	347.107	0.7183	99.79	10	0.9952	79.29	988
Erlotinib	394.177	0.7837	99.85	7	0.9775	61.04	1859
Sunitinib	399.220	0.6845	99.73	13	0.9921	72.23	1325
Pazopanib	438.171	0.8853	99.98	1	0.9837	64.60	1689
Gefitinib	447.160	0.8362	99.92	4	0.9948	78.22	1039
Sorafenib	465.094	0.8328	99.90	5	0.9951	79.04	1000
Dasatinib	488.164	0.6400	99.62	18	0.9980	92.10	377
Imatinib	494.267	0.7611	99.81	9	0.9766	60.64	1878
Dabrafenib	520.109	0.5499	97.78	106	0.9964	83.29	797
Lapatinib	581.143	0.6715	99.69	15	0.9775	60.97	1862
Trametinib	616.086	0.1696	70.72	1397	0.9038	53.07	2239

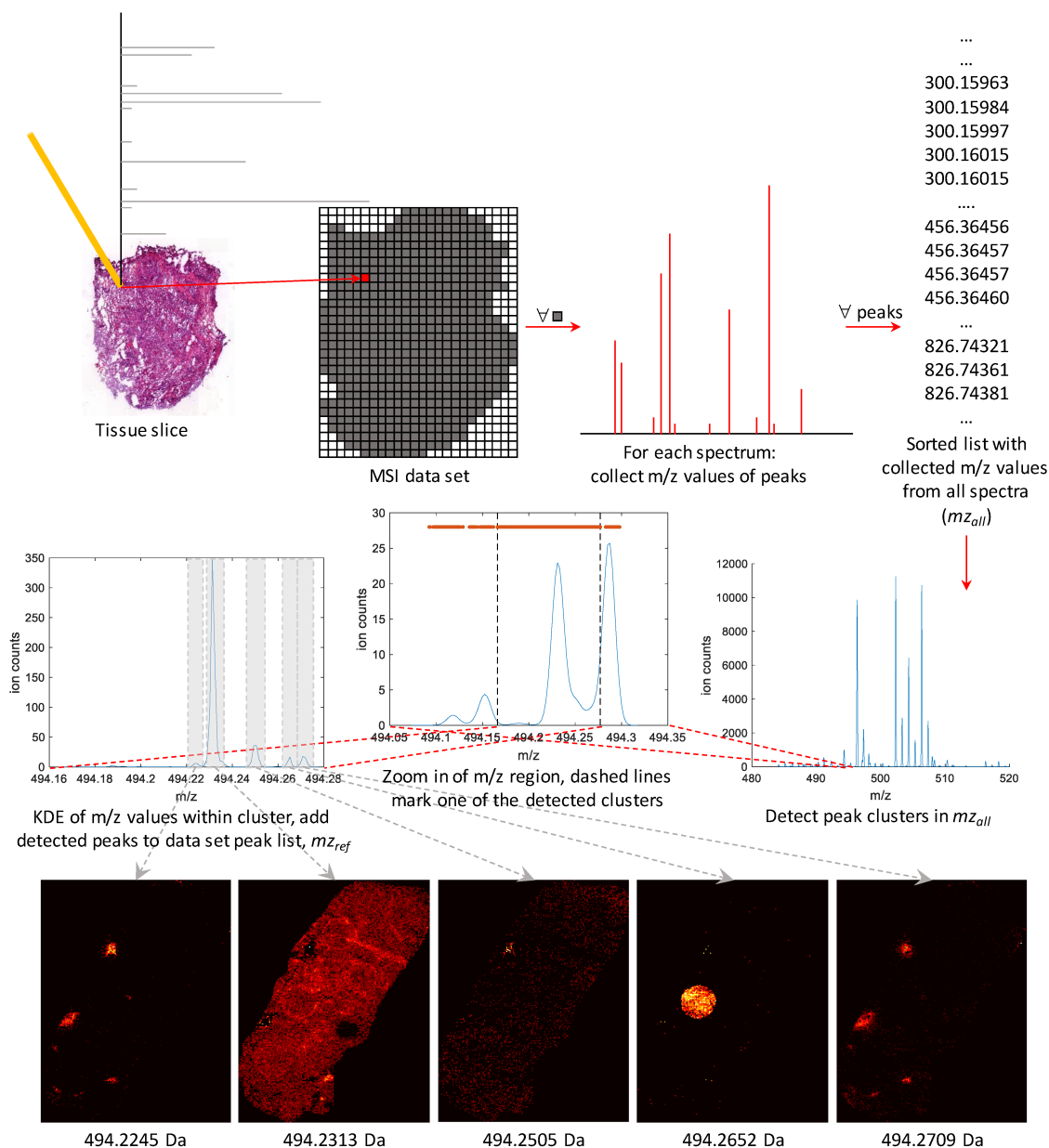
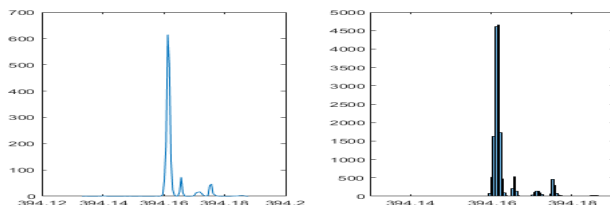
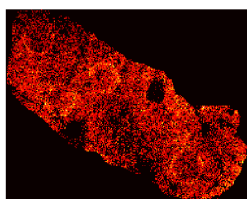


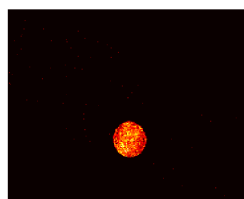
Figure 1: Flowchart of our peak picking algorithm. m/z values of peaks from each individual spectrum are collected and sorted in mz_{all} . We then identify clusters in mz_{all} as connected components in a directional graph. For each cluster we fit an optimized KDE to the distribution of m/z values. Data set peaks are obtained as local maxima on the resulting KDE curve. Finally, the level of structure in the ion images corresponding to the data set peaks is estimated and used to filter out noise peaks. The peak corresponding to the center ion image, at $m/z = 494.2505$ is an example of one filtered out in the last step.



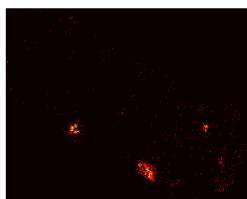
(a) **Left:** Optimized KDE curve. **Right:** Histogram of m/z distribution.



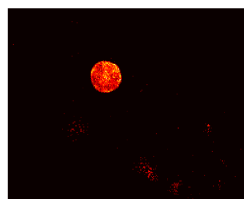
(b) 394.161 m/z



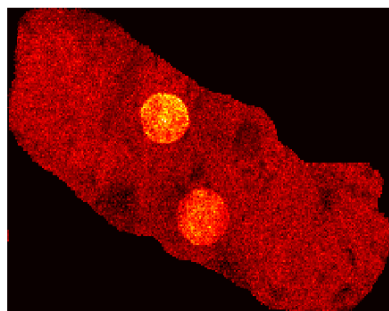
(c) 394.166 m/z



(d) 394.172 m/z



(e) 394.176 m/z



(f) Ion-image obtained by binning spectra between 394.15 and 394.20 m/z .

Figure 2: **(2a)** the distribution of m/z peak values within the cluster containing erlotinib (m/z 394.176). **(2b - 2e)** the ion images that correspond to the four peaks on the KDE curve. **(2f)** the ion image obtained by binning the spectra between 394.15 and 394.20 m/z ; this image demonstrates how four signals can be mixed in the same ion image and even when a relatively narrow m/z window is used.

Table 2: The VE and SC scores of the ion images corresponding to the eleven compounds reported by Römpp et al.¹⁰ and their corresponding rank among the 1 053 candidate ion images that remain after removing those with fewer than 200 nonzero pixels.

Compound	Mass	VE	Percentile	Rank (VE)	SC	Percentile	Rank (SC)
LPC (16:0), $[M + K]^+$	535.296	0.1770	92.76	74	0.9897	94.52	56
LPC (18:0), $[M + K]^+$	562.327	0.2732	98.14	19	0.9964	99.12	9
heme b, M^+ ,	616.177	0.2385	96.67	34	0.9261	70.84	298
unknown	713.452	0.0911	75.93	246	0.9444	73.68	269
SM (16:0)	742.531	0.2140	95.50	46	0.9953	98.24	18
unknown	743.548	0.1921	94.42	57	0.9691	84.34	160
PC(32:1), $[M + K]$	770.507	0.2688	97.95	21	0.9814	88.85	114
SM(18:0), $[M + K]$	770.565	0.1439	87.87	124	0.9849	90.90	93
PC (32:0), $[M + K]^+$	772.525	0.3177	98.83	12	0.9975	99.80	2
PC (34:1), $[M + K]^+$	798.541	0.3383	99.02	10	0.9979	99.90	1
PE(38:1)	812.557	0.1623	91.39	88	0.9909	95.21	49

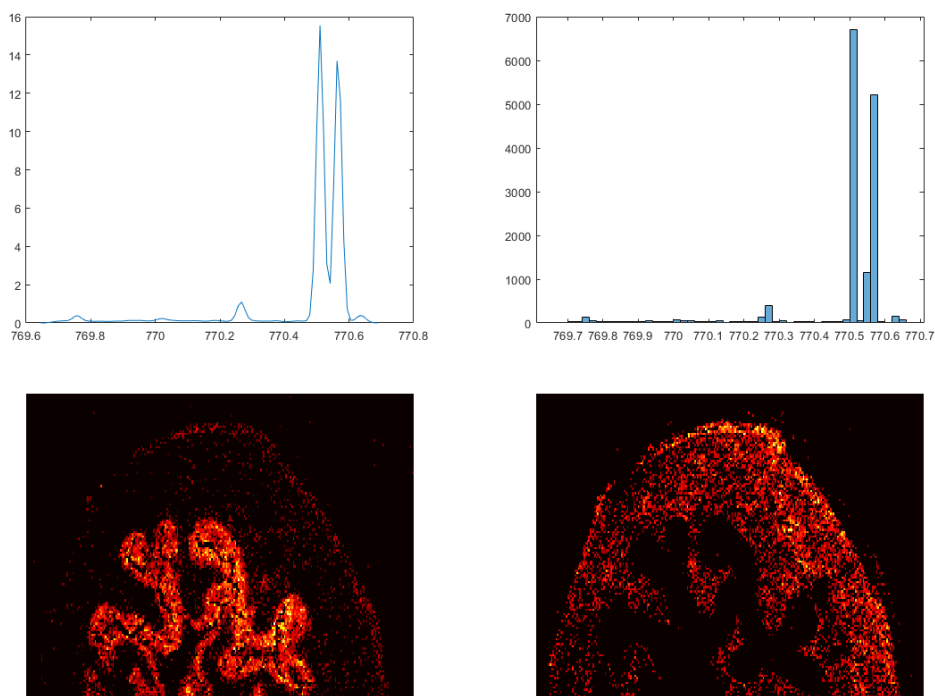


Figure 3: Distribution of peak m/z values within the cluster containing PC(32:1) (770.5109 m/z) and SM(18:0) (770.5609 m/z). The ion images corresponding to the two highest peaks on the KDE curve are shown in the bottom left and bottom right.

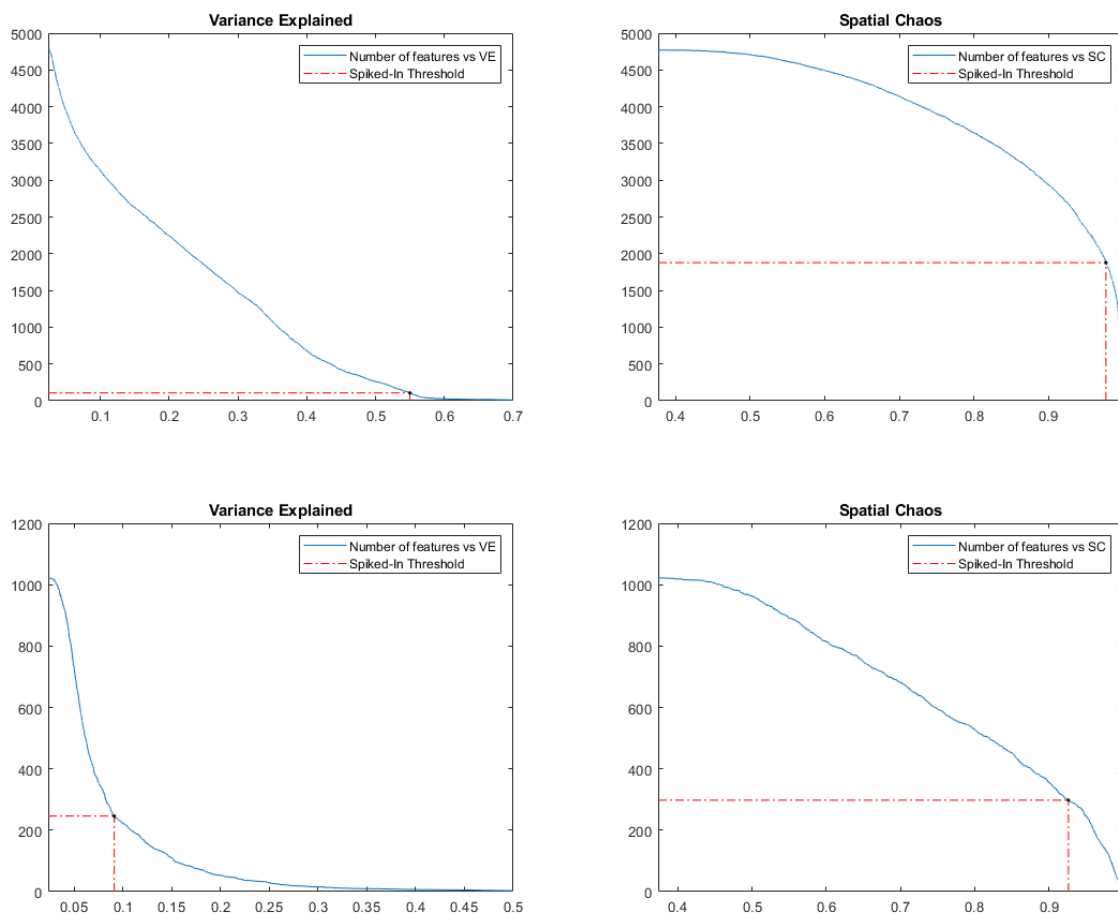


Figure 4: Number of ion images surviving varying thresholds on the VE and SC scores in the two data sets. Dashed lines mark the lowest scores (excluding the low quality image for m/z 616.127) of the ion images corresponding to the drugs in the spiked data set (**top**) and known compounds in the mouse bladder data set (**bottom**).

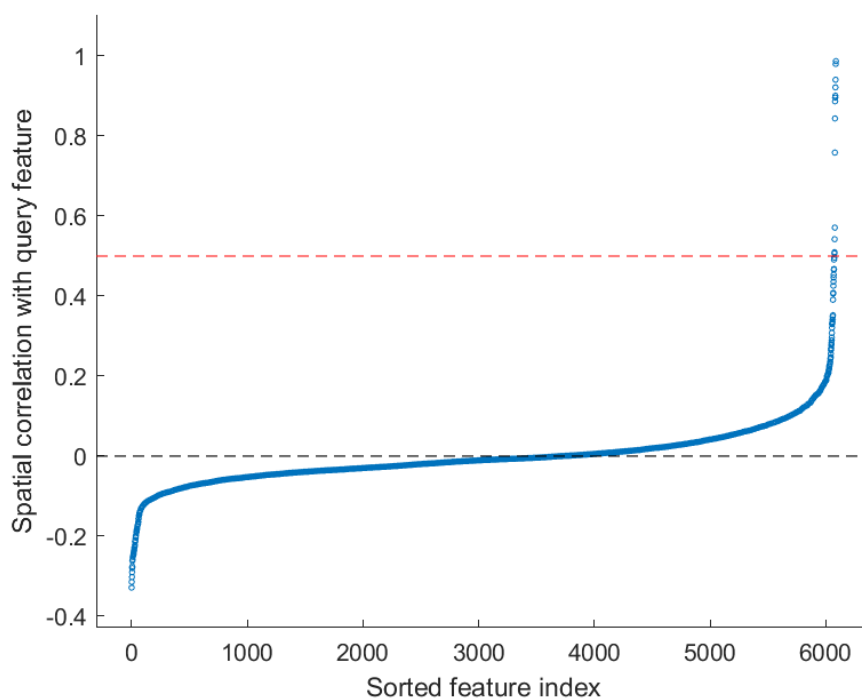
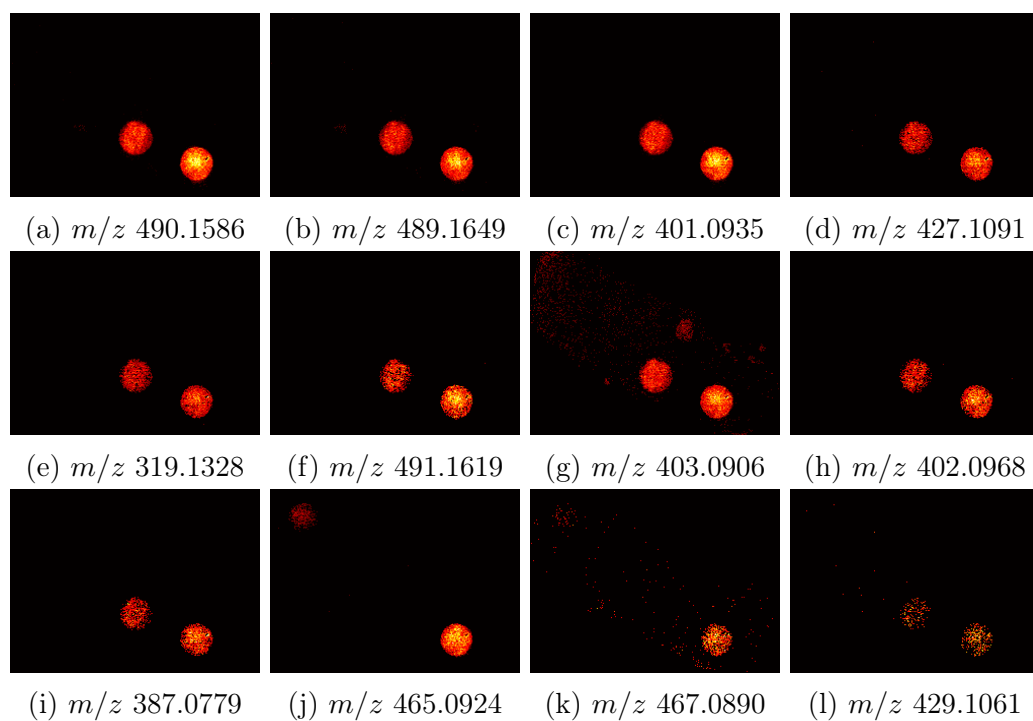


Figure 5: **top:** The ion images of the 12 most correlated peaks to dasatinib's monoisotopic peak. **(a - i)** and **(l)** are isotopes or fragments of dasatinib while **(j)** and **(k)** are related to sorafenib. **bottom:** sorted Pearson correlation between all ion images and that of the monoisotopic peak of dasatinib.

Acknowledgement

We thank Frank Suits for his support and insightful discussions throughout the project.

Supporting Information Available

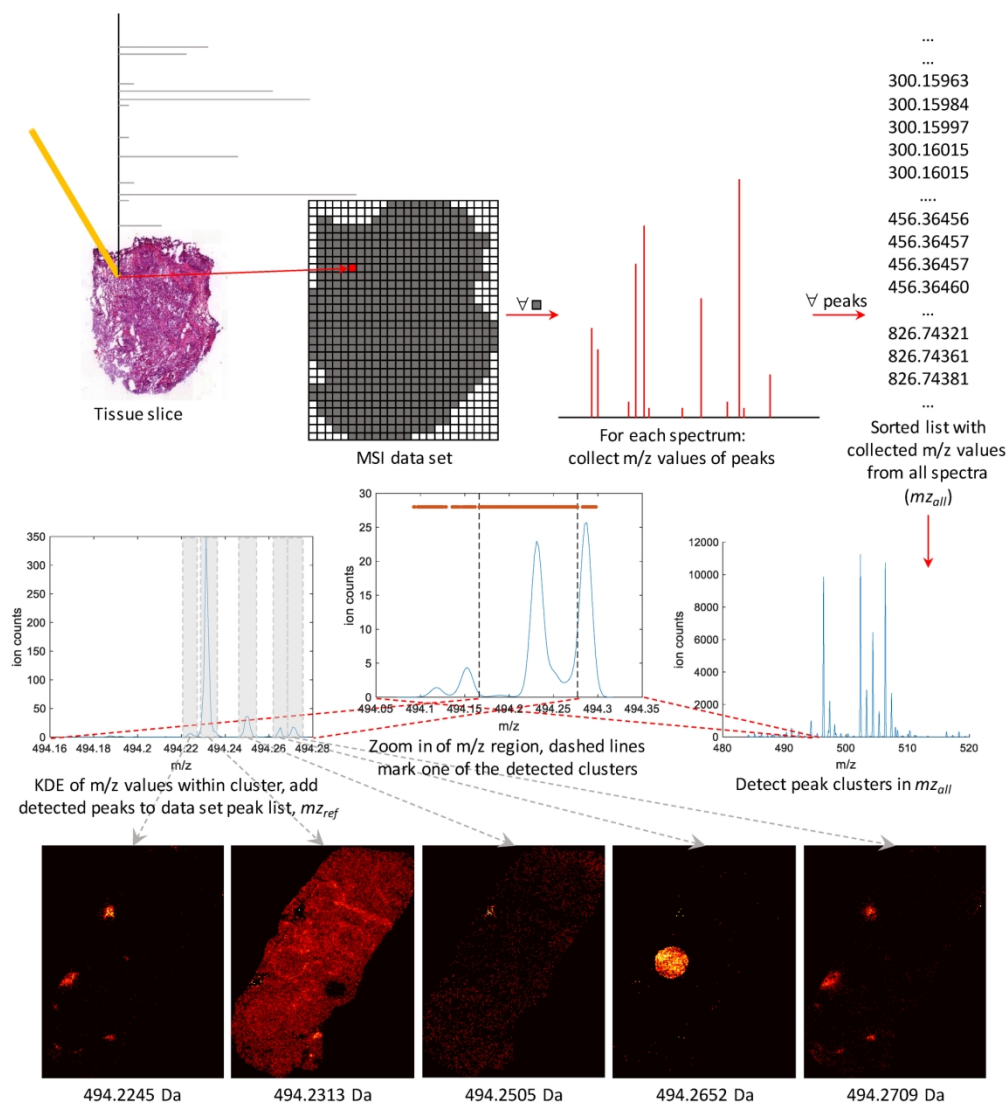
The following files are available free of charge as supporting information.

- Eriksson_et_al_supporting_information.pdf: methods and figures in Supporting Information.
- spikedin_drug_information.xlsx: table of correlating peaks for each spiked-in compound with structures and annotations (isotopes, fragments).
- drug_structures.pdf: structures of spiked-in drugs.

References

- (1) Jones, E. A.; Deininger, S.-O.; Hogendoorn, P. C.; Deelder, A. M.; McDonnell, L. A. *Journal of proteomics* **2012**, *75*, 4962–4989.
- (2) Gessel, M. M.; Norris, J. L.; Caprioli, R. M. *Journal of proteomics* **2014**, *107*, 71–82.
- (3) Bemis, K. D.; Harry, A.; Eberlin, L. S.; Ferreira, C.; van de Ven, S. M.; Mallick, P.; Stolowitz, M.; Vitek, O. *Bioinformatics* **2015**, *31*, 2418–2420.
- (4) Gibb, S.; Strimmer, K. *Bioinformatics* **2012**, *28*, 2270–2271.
- (5) Suits, F.; Fehniger, T. E.; Végvári, Á.; Marko-Varga, G.; Horvatovich, P. *Analytical chemistry* **2013**, *85*, 4398–4404.
- (6) Alexandrov, T.; Bartels, A. *Bioinformatics* **2013**, *29*, 2335–2342.

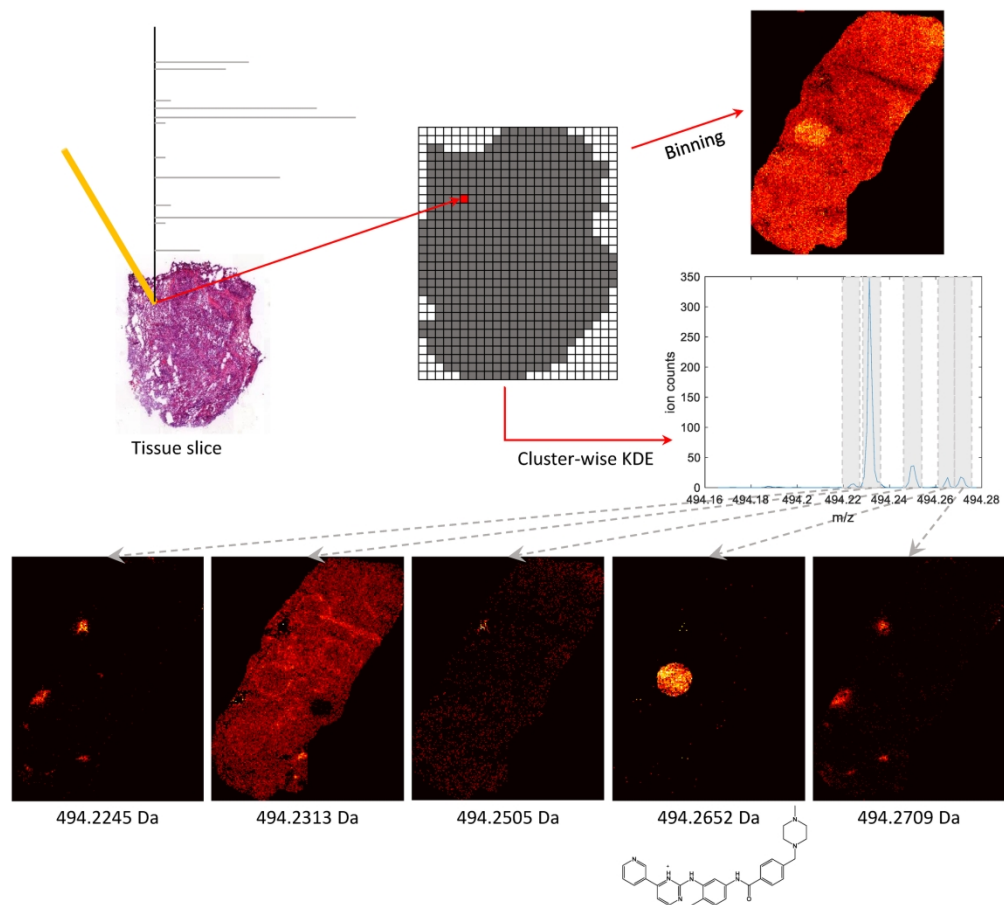
- 1
2
3 (7) Wijetunge, C. D.; Saeed, I.; Boughton, B. A.; Spraggins, J. M.; Caprioli, R. M.;
4 Bacic, A.; Roessner, U.; Halgamuge, S. K. *Bioinformatics* **2015**, *31*, 3198–3206.
5
6
7
8 (8) Palmer, A.; Phapale, P.; Chernyavsky, I.; Lavigne, R.; Fay, D.; Tarasov, A.; Kovalev, V.;
9 Fuchser, J.; Nikolenko, S.; Pineau, C.; Becker, M.; Alexandrov, T. *Nature methods*
10 **2017**, *14*, 57.
11
12
13
14 (9) Inglese, P.; Correia, G.; Takats, Z.; Nicholson, J. K.; Glen, R. C. *Bioinformatics* **2018**,
15 *35*, 178–180.
16
17
18
19 (10) Römpp, A.; Guenther, S.; Schober, Y.; Schulz, O.; Takats, Z.; Kummer, W.; Spen-
20 gler, B. *Angewandte chemie international edition* **2010**, *49*, 3834–3838.
21
22
23
24 (11) Schramm, T.; Hester, A.; Klinkert, I.; Both, J.-P.; Heeren, R. M.; Brunelle, A.;
25 Laprévotte, O.; Desbenoit, N.; Robbe, M.-F.; Stoeckli, M.; Spengler, B.; Römpp, A.
26 *Journal of proteomics* **2012**, *75*, 5106–5110.
27
28
29
30
31 (12) Hoffman, E. D.; Stroobant, V. *West Sussex: John Wiley & Sons, Bruxellas, Bélgica*
32 **2007**, *1*, 85.
33
34
35
36 (13) Suits, F.; Hoekman, B.; Rosenling, T.; Bischoff, R.; Horvatovich, P. *Analytical chemistry*
37 **2011**, *83*, 7786–7794.
38
39
40
41 (14) Bowman, A. W.; Azzalini, A. *Applied smoothing techniques for data analysis: the kernel*
42 *approach with S-Plus illustrations*; OUP Oxford, 1997; Vol. 18.
43
44
45
46 (15) Fonville, J. M.; Carter, C.; Cloarec, O.; Nicholson, J. K.; Lindon, J. C.; Bunch, J.;
47 Holmes, E. *Analytical chemistry* **2012**, *84*, 1310–1319.
48
49
50
51 (16) Nemes, P.; Woods, A. S.; Vertes, A. *Analytical chemistry* **2010**, *82*, 982–988.
52
53
54 (17) Fehniger, T. E.; Suits, F.; Végvári, Á.; Horvatovich, P.; Foster, M.; Marko-Varga, G.
55 *Proteomics* **2014**, *14*, 862–871.
56
57
58
59
60



42 Figure 1: Flowchart of our peak picking algorithm. m/z values of peaks from each individual spectrum are
43 collected and sorted in mz_{all} . We then identify clusters in mz_{all} as connected components in a directional
44 graph. For each cluster we fit an optimized KDE to the distribution of m/z values. Data set peaks are
45 obtained as local maxima on the resulting KDE curve. Finally, the level of structure in the ion images
46 corresponding to the data set peaks is estimated and used to filter out noise peaks. The peak corresponding
47 to the center ion image, at $m/z = 494:2505$ is an example of one filtered out in the last step.

48 172x190mm (300 x 300 DPI)

60



TOC figure

211x190mm (300 x 300 DPI)